

© 2016 by Narender Gupta. All rights reserved.

AMERICAN GRADUATE ADMISSIONS:  
BOTH SIDES OF THE TABLE

BY

NARENDER GUPTA

THESIS

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Computer Science  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2016

Urbana, Illinois

Advisor:

Professor Dan Roth

# Abstract

This is a comprehensive study of graduate admission process in American universities. There are multiple entities involved in the process, out of which the most significant ones are:

- The candidate applying for admission in a department in a school
- The decision-makers acting upon the candidates' application

The goal of this study is to understand the admission process from each of these entities' perspective, and provide them decision-support models for their respective tasks. Although both of the entities interact through a common set of datapoints, i.e. candidate admission application, each of them works towards a very different goal. The juxtaposition of these two tasks provides a very interesting challenge which is hard to resolve deterministically. Solution to such a problem requires learning techniques which can find patterns, adapt according to the dynamic nature of problem, and produce results in a probabilistic fashion.

We study and model the graduate admission process from a machine learning perspective based on analysis of large amounts of data. The analysis considers factors such as standardized test scores, and GPA, as well as world knowledge such as university *similarity*, *reputation*, and *constraints*. Based on the targeted entity, learning problem is formulated as classification problem or ranking problem. During learning and inference, not only those features are considered which are available from the data directly, but also the hidden features which need to be incorporated generatively. Our experimental study reveals some key factors in the decision process and, consequently, allows us to propose a recommendation algorithm that provides applicants the ability to make an informed decision regarding where to apply, as well as guides the decision-makers towards a more efficient process.

*To my parents, Ram Pratap and Saroj.*

# Acknowledgments

To begin, I'd like to thank my advisor, Professor Dan Roth, for believing in me and guiding me through the moments of confusion. He not only enabled me to cross the technical hurdles, but also helped me understand the faculty's perspective towards this problem which was indispensable.

The next most significant person throughout this research was Aman Sawhney who ignited this idea, and helped me nurture it at every step. The fact that he got into his master's degree because of this work gives me immense satisfaction.

I want to thank Viveka Kudaligama for sharing her invaluable expertise on the graduate admissions, and helping me jump through administrative hoops without faltering.

Thanks to the Department of Computer Science at University of Illinois at Urbana-Champaign for offering me the Research Assistant position, providing me with the financial means to complete this project.

And finally, thanks to my parents, and numerous friends who endured this long process with me, always offering support and love.

# Table of Contents

<b>List of Tables</b> . . . . .	<b>vii</b>
<b>List of Figures</b> . . . . .	<b>viii</b>
<b>List of Abbreviations</b> . . . . .	<b>ix</b>
<b>Chapter 1 Introduction</b> . . . . .	<b>1</b>
<b>Chapter 2 Admission Process and Data</b> . . . . .	<b>3</b>
2.1 Graduate Admission Process . . . . .	3
2.2 Application Graph . . . . .	3
<b>Chapter 3 Related Works</b> . . . . .	<b>5</b>
3.1 Qualitative Methods . . . . .	5
3.2 Statistical Approaches . . . . .	5
<b>Chapter 4 Two Sides of the Table</b> . . . . .	<b>7</b>
4.1 Two Sides . . . . .	7
4.2 Solution . . . . .	8
4.2.1 Classification Problem . . . . .	9
4.2.2 Ranking Problem . . . . .	9
<b>Chapter 5 Student Perspective</b> . . . . .	<b>10</b>
5.1 Dataset . . . . .	10
5.1.1 Data Cleaning . . . . .	10
5.2 Supervised Learning . . . . .	12
5.2.1 Features: . . . . .	12
5.2.2 Ensemble Learning: . . . . .	12
5.3 Latent Variable Based Approach . . . . .	12
5.4 Experimental Evaluation . . . . .	15
5.4.1 Evaluation Metric . . . . .	15
5.4.2 Discriminative Classifiers . . . . .	17
5.4.3 Feature Selection . . . . .	20
5.4.4 Latent Variable Based Approach . . . . .	21
5.4.5 Understanding Institution Rankings . . . . .	23
5.4.6 Impact of Change in Application Year . . . . .	26
5.4.7 Which Universities Go Together . . . . .	27
5.5 Recommendations . . . . .	28
<b>Chapter 6 University Perspective</b> . . . . .	<b>30</b>
<b>Chapter 7 Analysis and Discussion</b> . . . . .	<b>31</b>

Chapter 8 Conclusion and Future Work . . . . .	33
References . . . . .	34

# List of Tables

5.1	Statistics about Edulix Data . . . . .	11
5.2	Classification Context for F1 calculation . . . . .	17
5.3	Binary Classifier parameters tuned by grid search . . . . .	17
5.4	F1 over different classifiers and schemes . . . . .	20
5.5	Discriminative Power of each feature . . . . .	20
5.6	F1 without each feature. Less F1 due to missing feature indicates more discriminative power of that feature. . . . .	23
5.7	Gain in F1 score due to EM clustering . . . . .	23
5.8	Gain in $F1_{admit}$ due to various rank-lists (Average over 100 iterations) (In the order of magnitude of $10^{-3}$ ) . . . . .	26
5.9	Interesting similar universities based on Kulczynski score . . . . .	28
5.10	Universities that go together based on Apriori algorithm . . . . .	28



# List of Figures

2.1	An oversimplified view of ideal data for the graduate admissions problem. . . . .	4
4.1	Difference in problem understanding of students and universities. . . . .	8
5.1	$F1_{sparse}$ as a function of presence of sparse label. The curve is fitted with moving average function of window size 5. . . . .	18
5.2	$F1_{admit}$ as a function of Acceptance Ratio. The curve is fitted with moving average function of window size 5. . . . .	19
5.3	$F1_{admit}$ as a function of Graduate University US News Rank. The curve is fitted with moving average function of window size 5. . . . .	21
5.4	$F1_{admit}$ as we keep on increasing features on top of GPA. Refer to Table 5.5 for feature corresponding to index number on X-axis. . . . .	22
5.5	Gain in F1 due to various rank-lists . . . . .	25

# List of Abbreviations

ASU	Arizona State University
AWA	Analytical Writing Analysis
CalTech	California Institute of Technology
CMU	Carnegie Mellon University
CS	Computer Science
CSRL	Consciously Shuffled Rank List
CSU	Chicago State University
GMU	George Mason University
GPA	Grade Point Average
GRE	Graduate Record Examination
HCI	Human Computer Interaction
IELTS	International English Language Testing System
IU Bloomington	Indiana University, Bloomington
MCS	Master of Computer Science (Professional Master's degree program at UIUC)
ORL	Original Rank List
RSSDRL	Randomly Shuffled Same Distribution Rank List
RSUDRL	Randomly Shuffled Uniform Distribution Rank List
SJSU	San Jose State University
SUNY Bingham	State University of New York, Binghamton
SUNY Buffalo	State University of New York, Buffalo
SUNY Stony	State University of New York, Stony Brook
TAMU	Texas A & M University, College Station
TOEFL	Test Of English as Foreign Language
UC Boulder	University of Colorado, Boulder
UChicago	University of Chicago

UCLA	University of California, Los Angeles
UCR	University of California, Riverside
UCSC	University of California, Santa Cruz
UG	Undergraduate
UIC	University of Illinois at Chicago
UIUC	University of Illinois at Urbana-Champaign
UM Twin	University of Minnesota Twin Cities
UMD	University of Maryland, College Park
UN Vegas	University of Nevada, Las Vegas
UNCC	University of North Carolina, Charlotte
URI	University of Rhode Island
UT Austin	University of Texas, Austin
UWisc	University of Wisconsin, Madison

# Chapter 1

## Introduction

According to National Center for Science and Engineering Statistics from National Science Foundation, more than 600,000 graduate students were enrolled in graduate programs in America in 2014, in science and engineering fields alone [National Science Foundation, 2014]. Adding medical and other disciplines easily crosses the figure of a million. Every year, hundred of thousands of these students, and many others who don't get admission, apply to American graduate programs and, in the process, discover that there is a dearth of reliable sources to aid them in making an informed decision. There are several sources that provide admission related statistics, but do not cater to individual needs, thus, leaving the applicant with the only option of guessing and hoping for the best [US News, 2015, QS Quacquarelli Symonds Limited, 2015]. An equivalently difficult, but different, problem is faced by the members of admission committee at each school who receive thousands of candidate applications every year to evaluate. These evaluators, based on rules or their experiences, analyze every application and recommend weather to *Admit* or *Reject* the candidate. Despite hundreds of years of existence of the process, it takes several weeks at the least, and many a times several months, for most of the schools to complete this evaluation. Even though learning models have been used in variety of real-world applications, there is surprisingly little literature available on understanding admission dynamics and decision making. This research attempts to fill the gap by providing a decision support-model to both of the involved entities i.e. candidates and evaluators. These models enable them to make informed choices by taking into account several factors such as university acceptance rate, similarity, constraints, as well as admission trends.

Rest of this work is organized in several chapters where each chapter talks about a semantically coherent set of ideas, experiments or observations. Chapter 2 explains admission process and several variables involved in the process. It also talks about the kind of data one needs to solve this problem. Chapter 3 goes through the attempts researchers have made to solve this problem and their respective drawbacks. In Chapter 4, we present the idea that different entities involved in the process have different problems, and, hence, each require their own customized solutions. Chapter 5 deals with the problem from a student's perspective. It defines the problem mathematically, delineates several experiments conducted on the data, and presents the

findings. In this chapter, we also propose a recommendation algorithm for students to use so that they can make an informed decision while choosing graduate schools. Chapter 6 gauges the problem from evaluators' perspective, again involving mathematical definition, experimental observations and results. Chapter 7 analyzes the problem further with the intention of revealing interesting patterns, and discusses the findings of previous chapters in a coherent fashion. In chapter 8, we conclude this work by connecting all the dots, and discussing where the future research should or could go.

This thesis will focus on student's perspective of the problem. The work was done on perspective of the university, as part of being employed by the University of Illinois at Urbana-Champaign, but will not be part of this thesis.

## Chapter 2

# Admission Process and Data

### 2.1 Graduate Admission Process

A typical university application packet comprises of transcripts, standardized test scores, letters of recommendation, a statement of purpose that expresses student's aims, ambitions and research interests, and descriptive answers to a few additional questions. Test scores include GRE, language test scores - such as TOEFL or IELTS etc. Universities, then, evaluate application packets based on confidential rules or heuristics and release decisions. Since requirements, deadlines and the specific process to meet them is university specific, the applicant needs to first choose the universities he would apply to.

Given the uncertainty, a naive solution is to apply to a large number of universities. But, the more the number of applications, the higher the investment of time and energy. This also implies a large monetary investment, which is a major concern for applicants from developing countries. One of the strategies to circumvent this is to categorize the universities into buckets so that one only applies to a few representatives from each category. A popular scheme includes three categories: *Dream*: where the chances of admission are slim; *Reach*: where the chances of admission are decent; *Safety*: where there is a fair certainty of being accepted [Krishnamoorthy, 2013]. Multiple admission offers resulting from these decisions allow the applicant to choose, suboptimally, their best option. The description of these categories is very subjective, and even more subjective is the applicant's ability to predict his probability of *Admit* i.e. chances of being admitted to a given program [Krishnamoorthy, 2013]. This prediction is generally based on hearsay or semi-informed opinions, resulting in confusion and a waste of resources for both the applicant as well as the university.

### 2.2 Application Graph

Since a candidate applies to many universities for admission, and each university gets applications from hundreds, or even thousands, of applicants, an ideal dataset would be a bipartite graph between the nodes of types - candidates and universities. We call this graph as the *Application Graph*. Each of the nodes

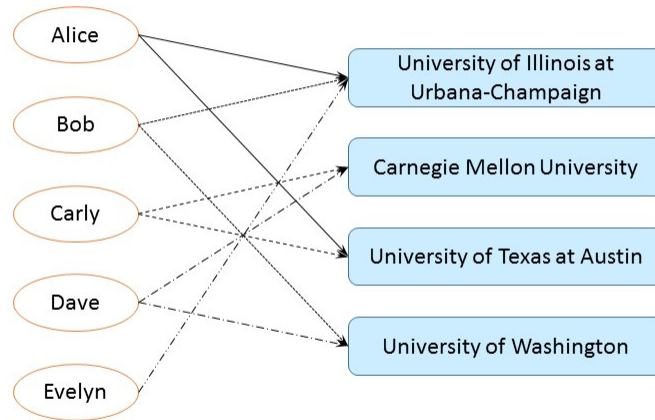


Figure 2.1: An oversimplified view of ideal data for the graduate admissions problem.

has some associated metadata with it. For a candidate, it includes all of the parts of the application discussed above i.e. transcripts, test scores, letters of recommendation, statement of purpose and so on. For a university, the metadata could include the number of programs of study offered, differences in nature of these programs, number of students admitted in each program, financial constraints, geographical or political constraints, and so on. An oversimplified view of this dataset would look something like Fig 2.1.

# Chapter 3

## Related Works

### 3.1 Qualitative Methods

Some researchers have briefly reflected on the process of decision making from university's perspective, but only qualitatively [Raghunathan, 2010], [Posselt, 2016]. Raghunathan talks about his experience as part of the admission committee in Computer Science department during his graduate studies at Stanford. He documents several factors considered during the decision-making and provides his opinion regarding their importance towards an *Admit*. Posselt's study provides more coverage in terms of number of universities, and the number of departments per university. Posselt focuses only on doctorate (PhD) admissions while Raghunathan had talked about only Master of Science (MS) applications. Also, Posselt is able to garner faculty opinion through surveys, interviews and direct observation of admission process, whereas Raghunathan's opinion is that of a graduate student and might be limited because of restrained participation. While both of these studies seem complimentary, none of them provides concrete data points about critical questions pertinent to the domain. Findings in these works are assimilation of individual opinions and fall short of any statistical measure.

### 3.2 Statistical Approaches

To overcome these limitations, certain researchers have tried to look at the problem from a statistical point of view. The work done by Waters et al. models the problem from the university's perspective quantitatively [Waters and Miikkulainen, 2013]. They used a learning approach to aid the university admission committee by identifying the candidates that are unlikely to be offered admission. Their model is quite simplistic, considering the problem as a straightforward classification problem (via Logistic Regression) without attempting to reveal the diverse and rich patterns in the data.

Works such as Bruggink et al. and Moore et al. utilize domain knowledge to build statistical models [Bruggink and Gambhir, 1996, Moore, 1998]. Bruggink et al. model undergraduate university admissions



to a private liberal arts college. The model treats application components as independent variables and assumes the decision to be dependent (*Admit* or *Reject*) on these. The independent variables include GPA, SAT scores, other academic scores and extracurricular factors all of which have been quantized. Beyond strong statistical assumption, this approach assumes that the modelled university (and its application pool) provide a good representation of the whole distribution. Our study shows that this is not the case because different universities focus on different features, and hence produce decisions differently. Moore et al. model the problem with rule induction using ID3 algorithm. Such an approach without care for bounded depth is prone to overfitting. Similar to Bruggink et al. the model is centered around one university and considers very small applicant sample size.

More importantly, all of the above approaches are university centric and do not provide any support to the decision process of applicants. The primary reason for the lack of such endeavours is the unavailability of a relevant dataset. One contribution of our work is the creation of such a dataset, which we will make available to the community. The dataset allowed us not only to determine the acceptability of an application but also suggest better choices. We conclude that learning a decision model should be done separately for each university, if possible.

# Chapter 4

## Two Sides of the Table

### 4.1 Two Sides

While the trend emerges from previous studies that modeling the problem as a regression problem produces some results, assuming so without due analysis would be naïve. Since there are primarily two types of entities involved in the graduate admission process, applicants and universities, it is imperative that we study the problem from each of their point of views. An applicant wants to get *Admit* from the best possible university, while a university wants to *Admit* the best possible candidates that it can. Both of these entities have access to different resources, have different understanding of the problem, and are working towards very different goals.

An applicant has access to his entire professional data, personal interests and biases, but can only see publicly accessible information related to universities. Similarly, a university has access to all of its internal data, previous years' students' data, requirement specifications and resources, but can only access the part of candidate's data provided in the application packet. A student is constrained by either financial limitations (each application has a fee associated with it), or personal biases such as geographical preferences etc. A typical university has relatively more resources in terms of financial wealth, but is limited in terms of the **Capacity** (number of candidates it can *Admit*), or **Human Effort** (admission committee can process a limited number of applications because of the manual process) etc.

Figure 4.1 shows a basic scenario explaining difference in problem understanding of each of these entities. Alice, Bob, Carly, Dave, and Evelyn (assumed names for the purpose of illustration) are the candidates who applied to a few universities. Each arrow indicates an application which connects the applicant to the university. From an applicant's perspective, all of the admissions are independent of each other. Since Alice applied to UIUC and UT Austin, and these universities, don't talk to each other regarding Alice's admission decision, Alice's assumption holds true. At the same time, a university's perspective is that all of the applications are dependent on each other. Suppose UIUC got applications from Alice, Bob and Evelyn, but can offer *Admit* to only 1 applicant because of capacity limit. In such a case, no matter how good an

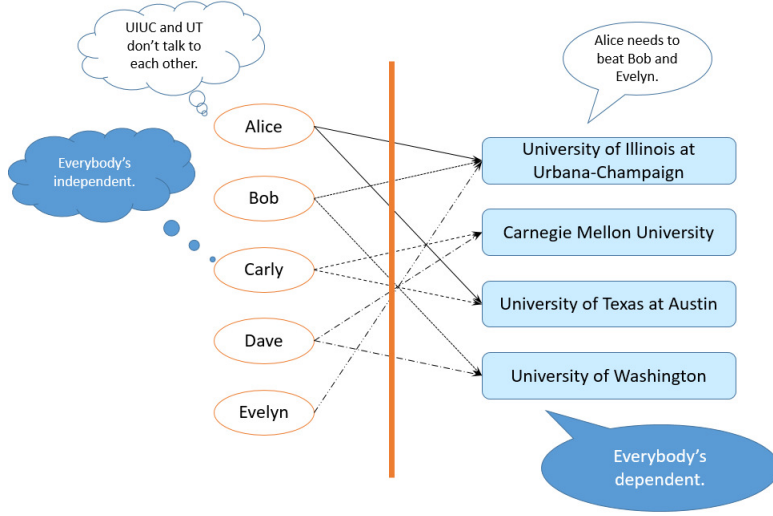


Figure 4.1: Difference in problem understanding of students and universities.

application is in its own standing, it won't get an *Admit* unless it outperforms the other two applicants. Since each application decision is relative to the peer applications, university's assumption also holds true.

So, we have two entities who are involved in the same process, but are working for different goals. This, and difference in access to resources, cause each of them to perceive the process differently, and hence act in distinct manners.

## 4.2 Solution

In an ideal world, if one had access to the complete Application Graph, it would be possible to model the complete process as a joint graphical model where each of the entity's constraints would control the interaction between itself and its neighbours. In a practical world, however, there are more than one reasons which make this data collection task next to impossible. There is no central system which keeps track of all of the applications flowing into the system via different universities. No two universities collaborate to align the received applications. And an application might contain sensitive information related to the candidate such as gender, race or ethnicity, making it harder for the universities to release the application data in public domain. Due to these reasons, construction of even a partial Application Graph is a grueling task.

Despite the lack of Application Graph, we can collect information related to applications from the perspective of each of the involved entities. Based on each of their perspective, we model the admission process as following two problems:

### 4.2.1 Classification Problem

Admission candidates share parts of their application information on the online website portals such as [The Grad Cafe, 2015], [Kassegne, 2016] etc. Information collected from such websites is an approximation of students' perspective. As discussed in Section 4.1, all of the admission decisions from students' perspective are independent of each other, the problem can be posed as a supervised binary classification problem with the labels  $\{Admit, Reject\}$ . All of the metadata from application can be used to extract features and then can be fed into a classification algorithm. The benefit of this modeling is that it uses the same independence assumption which is practically valid, and is flexible enough that it allows experimentation with multiple classification algorithms. The goal of this modeling is to identify as many correct decisions as possible so that the student has a fair estimate of his chances of getting into a school beforehand.

### 4.2.2 Ranking Problem

Universities, on the other hand, are working towards a different goal. They have richer data (because online portals have access to limited parts of application), but are interested in the top candidates based on their capacity. In such a case, the process can be modeled as a supervised ranking problem where *Admit* is ranked above *Reject*. Historical admission data can be used for training the ranking algorithm. The goal here is to optimize human effort so that the same admission process takes less time than a random application pool. In an ideal case, the ranking algorithm will rank all of the *True Admits* above all of the *True Rejects* so that the admission committee needs to look at only the top  $k$  applications where  $k$  is the capacity of university.

# Chapter 5

## Student Perspective

### 5.1 Dataset

There are several online resources where applicants share their admission experiences [Kassegne, 2016, The Grad Cafe, 2015]. We scraped the data present on Edulix [Kassegne, 2016]. It is an active resource which hosts applicant profiles from all over the world. GRE scores, undergraduate university name, GPA, TOEFL scores and other accomplishments such as work experience and research publications pertinent to the graduate admissions are reported in the profile. In addition, users mention the universities that they applied to and the result of each application (*Admit*, *Reject* or *Result Not Available*). Since the data is self reported, it had some erroneous records. We fixed these problematic entries using techniques explained in Section 5.1.1. In this paper we focus on modeling admission to computer science graduate programs, which form the plurality of our data. Our experiments are conducted only on this subset of the data. A few of the statistics of the resulting dataset for the computer science applications are in Table 5.1.

#### 5.1.1 Data Cleaning

Since the data is self reported, it had some erroneous records. We fixed these problematic entries using following techniques:

- *Missing Values*: If the record did not have mandatory fields such as GRE Verbal & Quantitative, or GPA, then such a record was completely deleted. Missing year and missing term were replaced with the mean of existing values.
- *Out of Range Errors*: Records which had invalid values such as negative GPA, or out of range value for GRE score components were removed in a rule-based fashion.
- *Illegal Data-types*: Records having incompatible data-types such as string for GPA, or numerical value for degree were removed.

Table 5.1: Statistics about Edulix Data

<b>General statistics</b>	
Total number of users before sanitization	36,207
Total number of users after sanitization	26,148
<b>Statistics for CS related dataset</b>	
Number of users	10,788
Application year range	[2001 2015]
Median Application Year	2013
Most Frequent Application Term	Fall
Number of universities with reported data	313
Number of applications per student (Mean)	6
Number of applications per university (Mean)	51
Number of undergraduate universities	2353
Degrees sought	[MS, PhD]

- *Quantitative Normalization*: Undergraduate institutions all over the world follow different scales for reporting GPA. Similarly, GRE and TOEFL tests have undergone various scale transformations over the years. As a standardization measure, we mapped these fields linearly to a scale of [0,100].
- *Categorical Normalization*: Sometimes students refer to the same degree using various acronyms such as *BE* or *BEng* for *Bachelor of Engineering*. Similar is the case for departmental majors such as *CS* or *Comp Science* for *Computer Science*. Also, sometimes there are spelling mistakes and typos in these fields. Each of these fields is, hence, normalized to manually created codes e.g. `{ba,bs,be,btech etc.}` instead of original strings. The code assignment was done using regular expressions or a manual mapping of string value to code.
- *University Normalization*: An undergraduate university might be referred to by the differing names due to reasons such as usage of a popular acronym or spelling errors. We mitigated this problem by mapping the university names to their unique website homepage using web search engine. Under this scheme, a program queried a search engine which returned a ranked list of relevant URLs. These URLs were then filtered based on heuristics (e.g. remove blog or social network URLs, prefer *.edu* domains and so on) and most relevant URL was assigned to the university. These assignments were then manually verified for consistency.
- *Labels*: There were no such labels present in dataset as *Waitlisted*, or *Conditionally Accepted*. We excluded any application that was classified as *Result Not Available* because it simply represents either a missing value or a pending decision.

## 5.2 Supervised Learning

Each university offers binary decision to applicant (*Admit* or *Reject*) and this decision is independent of the decision of other universities. Hence, the overall problem can be modeled as a set of individual binary classification problems. Supervised learning algorithms can be trained using features extracted from a labeled dataset and evaluated for performance.

### 5.2.1 Features:

The dataset contains several fields such as standardized test scores and academic history records. We extract several numerical features from these such as GRE test scores (AWA, Verbal & Quantitative), undergraduate GPA, language test scores (TOEFL), as well as categorical features such as program applied to (e.g. MS, PhD), term (e.g. Fall, Spring) etc. These features are used for learning Logistic Regression, Support Vector Machine and Random Forest. Since most work available in the literature is confined to approaches we have mentioned so far, we'll regard it as the baseline for any novelties that we propose. Table 5.5 lists all of the features extracted from the applications.

### 5.2.2 Ensemble Learning:

To improve the system performance, we use ensemble learning. Training decision trees without bounded depth is prone to overfitting [Quinlan, 1986], [Murthy and Salzberg, 1995]. But we can hope to generalize better if we use several limited-depth decision trees using partial data, and then feeding each decision as a feature into another regularized classifier. We create  $d$  such limited-depth decision tree classifiers where each classifier is trained on a bootstrapped sample from the original dataset [Efron and Tibshirani, 1986]. Bootstrapping allows us to have a tunable number of approximate feature representations instead of low number of exact features. Constrained by the variance-bias trade-off, the second classifier captures variance of data through multiple underlying decision trees while keeping a limit on its own variance by choosing simple models such as linear separators e.g. soft-margin support vector machine. Corresponding to decision of each such classifier, we get a feature for the next classifier.

## 5.3 Latent Variable Based Approach

For a few universities, we noticed that our discriminative classifiers do not perform as expected. We attribute it to the fact that some universities offer multiple degree programs that might target different kinds of applicants and have different admission criteria, e.g. professional master's versus thesis master's program

at UIUC. Certain other universities, such as CMU, offer specific programs at same degree level such as master's degree in Machine Learning versus HCI which fall under the purview of Computer Science. These distinctions, however, are not captured in the dataset, and are thus *hidden* from our models. This distribution causes difficulties for linear separators which assume that data is coming from single source. In case the data is coming from multiple sources, it is impossible to linearly separate the data using a single classifier. To accommodate these phenomena, we use an *Expectation-Maximization (EM)* algorithm that allows us to learn a model with a latent variable, capturing the missing information. The rest of this section, describes how we model the problem in this case.

EM is one of the very effective techniques for finding a *Maximum Likelihood Estimate* with *hidden variables* [Dempster et al., 1977]. For the sake of brevity, we are not going into details of *EM*. We will be using a setup similar to the one Grove et al. used in [Grove and Roth, 2001].

We assume a latent (hidden) variable called Program Type,  $z$ , which might carry one of multiple values. These values can hold different semantics based on university, e.g. different for UIUC and CMU, and, hence, we refer to them simply by the value assigned by the model to the hidden variable such as  $1, 2 \dots k$  etc. Once we learn the most likely model with the hidden variable  $z$  taking  $k$  values we can use it in two different ways. In a soft-boundary setting, each student can be from either Program Type with some probability. In a hard boundary setting, the most likely cluster (program) completely owns the student record, and we can learn individual linear classifiers for each cluster. A student belongs to cluster  $z = i$  if:

$$P(student|z = i) > P(student|z = j), \forall j \neq i \quad (5.1)$$

where

$$\sum_j P(student|z = j) = 1 \quad (5.2)$$

We assume that an applicant *belongs* to a program  $z \in \{1, 2, \dots, k\}$  with probability  $\alpha_r = p(z = r)$ . Given this program, feature  $x_i$  of the applicant  $x$  is generated independently by a Gaussian distribution with model parameters  $(\mu_i, \sigma_i)$ . Hence, we aim to split Gaussian mixture of data into individual models using *EM* setting.

Let us define a hidden variable,  $z \in \{1, 2, \dots, k\}$ . A student record (sample),  $x$ , consists of  $(n+1)$  features, The likelihood of sample is given by

$$P(x) = \sum_z p(x|z)p(z) \quad (5.3)$$



Incorporating each feature probability, the likelihood of the data sample can be expressed as:

$$P(x) = \sum_z \left( p(z) \prod_i p(x_i|z) \right) \quad (5.4)$$

Starting with an initial set of parameters  $\theta$ , the probability that a data point  $x^j = x^j = (x_0^j, x_1^j, \dots, x_n^j)$  comes from each of the  $k$  values of  $z$  is given by:

$$P_r^z = p(z = r|x^j) = \frac{p(z = r) \prod_i p(x_i^j|z = r)}{\sum_z \left( p(z) \prod_i p(x_i^j|z = r) \right)} \quad (5.5)$$

Let  $p_i^z = p(x_i|z = r)$ , then we can compute the expected log-likelihood as follows:

$$\begin{aligned} E(LL) &\equiv E \left( \sum_j \log P(x^j|\alpha_z, p_i^z) \right) \\ &\equiv \sum_j E \left( \log P(x^j|\alpha_z, p_i^z) \right) \\ &\equiv \sum_j \left( \sum_z P_z^j \cdot \log P(z, x^j|\alpha_z, p_i^z) \right) \\ &\equiv \sum_j \left( \sum_z P_z^j \cdot \log(\alpha_z \cdot \prod_i p(x_i^j|z)) \right) \end{aligned} \quad (5.6)$$

Assuming numerical features to be generated from a Gaussian distribution with parameters  $(\mu, \sigma)$ , above equation can be expanded as:

$$\begin{aligned} E &\equiv \sum_j \left( \sum_z P_z^j \cdot \log \left( \alpha_z \cdot \prod_i \frac{1}{\sigma_i^z \sqrt{2\pi}} \exp \left( -\frac{(x_i^j - \mu_i^z)^2}{2(\sigma_i^z)^2} \right) \right) \right) \\ &\equiv \sum_j \left( \sum_z P_z^j \cdot \left[ \log \alpha_z + \sum_i \left( -\log \sigma_i^z - \frac{(x_i^j - \mu_i^z)^2}{2(\sigma_i^z)^2} \right) \right] \right) \end{aligned} \quad (5.7)$$

Differentiating with respect to all parameters, we can find new values of  $\alpha_z$ ,  $\sigma_i^z$ , and  $\mu_i^z$ , for which the expected (log) likelihood receives an extremal value. Since,  $\alpha$  can only assume  $k - 1$  independent values because of  $k$  states of  $z$ , we have:

$$\alpha_k = 1 - \sum_{i \in [1 \dots k-1]} \alpha_i \quad (5.8)$$

Using this equation and differentiating Eq 5.7) partially with respect to model parameters, we get fol-

lowing update rules:

$$\alpha_z = \begin{cases} \alpha_k \frac{\sum_j P_z^j}{\sum_j P_k^j} & \text{if } z \neq k \\ 1 - \sum_{i \in [1 \dots k-1]} \alpha_i & \text{if } z = k \end{cases} \quad (5.9)$$

$$(\sigma_i^z)^2 = \frac{\sum_j \sum_z P_z^j (x_i^j - \mu_i^j)^2}{\sum_j \sum_z P_z^j} \quad (5.10)$$

$$\mu_i^z = \frac{\sum_j \sum_z P_z^j x_i^j}{\sum_j \sum_z P_z^j} \quad (5.11)$$

Using above update rules, and various initializations for the latent variables, we performed multiple experiments with EM. Current results are reported for  $z=2$  in a hard-boundary setting.

## 5.4 Experimental Evaluation

Our experimental study is designed to investigate the following issues:

- The ability to make reliable prediction on sparser label for a university for any student.
- The ability to make reliable prediction on whether a specific student can be admitted to a given program.
- Our ability to identify sub-programs in a given university, and its significance on the performance of our admission model.
- Understanding the contribution of factors to admission.
- Understanding the differences among universities in terms of their admission decisions.

### 5.4.1 Evaluation Metric

Before we evaluate performance, it is important to understand the appropriate performance metric for this task. Prediction accuracy is a biased metric in this case, because most universities have moderate to heavy imbalance in terms of the presence of labels (*Admit*, *Reject*) and a high accuracy doesn't necessarily indicate effective learning. Some of the universities have as low *Acceptance Ratio* as 8% while others accept more than half of the candidates that apply.

$$Acceptance\ Ratio = \frac{count(Admit)}{count(Admit) + count(Reject)}$$

In such an imbalanced label distribution, simple baseline of dense label assignment as prediction will result in high accuracy, without the need to learn anything. Hence, we choose F1 as our evaluation metric which takes into account not just the correct number of predictions made for the label (*Precision*), but also the ratio of predicted true labels out of total true labels (*Recall*). A natural choice for the label is *Admit*, since it is what any applicant cares about. But there are cases when it makes the problem suspiciously easier from theoretic standpoint. Consider two simple cases with following label distributions in data, and corresponding simple baselines where we blindly assign *Admit* prediction to each of examples:

1. *Admit : Reject* = 10 : 90

$$Precision = \frac{1}{10}, Recall = 1, \mathbf{F1(Admit)} = \frac{2}{11} = \mathbf{0.18}$$

2. *Admit : Reject* = 90 : 10

$$Precision = \frac{9}{10}, Recall = 1, \mathbf{F1(Admit)} = \frac{18}{19} = \mathbf{0.95}$$

It is evident that valuating F1 on the sparse (less frequent) label provides stricter bounds on the performance. Hence, from a machine learning perspective, it is easy and uninteresting to predict for the denser label, but from a student's perspective, it makes sense to predict *Admit*. Hence, we evaluate our models on both of the metrics i.e. F1 over sparser label ( $F1_{sparse}$ ), as well as F1 over *Admit* ( $F1_{admit}$ ).

Let  $U = \{1, 2, \dots, N\}$  is the set of  $N$  universities. Using Table 5.2,  $F1^{(i)}$  for the university  $i$  is defined as:

$$F1^{(i)} = \frac{2 \times Precision^{(i)} \times Recall^{(i)}}{Precision^{(i)} + Recall^{(i)}}$$

$$Precision^{(i)} = \frac{TP}{TP + FP}$$

$$Recall^{(i)} = \frac{TP}{TP + FN}$$

$F1$  over the set  $U$  is defined as the statistical mean of individual  $F1$ s i.e.

$$F1_{sparse} = \langle \{F1_{sparse}^{(i)} \forall i \in U\} \rangle$$

$$F1_{admit} = \langle \{F1_{admit}^{(i)} \forall i \in U\} \rangle$$

The rest of this section describes the details of our experimental study, and its results. In all our experiments we are reporting average F1 over 5-fold cross-validation.

Table 5.2: Classification Context for F1 calculation

		True Condition	
		Admit	Reject
		TP	FP
Predicted Condition	Admit	TP	FP
	Reject	FN	TN

### 5.4.2 Discriminative Classifiers

We ran multiple experiments with various classifiers such as - SVM with linear kernel or Radial Basis Function kernel [Chang and Lin, 2011], [Smola and Schölkopf, 2004], Logistic Regression, Adaboost with decision trees [Freund and Schapire, 1995] [Breiman et al., 1984], and Random Forest [Breiman, 2001]. We also experimented with all of the above classifiers with balanced class weights by adjusting them inversely proportional to the class frequency in input data. Each of these setups is used in two different settings:

- **Simple features:** Features extracted from student records
- **Tree features** Training  $d$  decision tree classifiers with bounded-depth and then using predictions of these classifiers as features. Each decision tree is trained on uniformly sampled 50% data restricted to maximum depth of 3.

Using grid search in the range [10,1000], we found that  $d=60$  yields maximum average F1 over all universities. These classifiers can provide the probability of the label as well which we utilize in Section 5.5. These experiments were conducted using Scikit-learn [Pedregosa et al., 2011]. Table 5.3 lists the tuned parameter values for different classifiers.

Table 5.3: Binary Classifier parameters tuned by grid search

Classifier (Param)	Range	Best Value
SVM(Linear): C	[1,100]	29.55
SVM <sub>balanced</sub> (Linear): C	[1,100]	26.29
SVM(RBF): C	[1,100]	3.68
SVM(RBF): $\gamma$	[0.01,1.0]	0.05
SVM <sub>balanced</sub> (RBF): C	[1,100]	2.59
SVM <sub>balanced</sub> (RBF): $\gamma$	[0.01,1.0]	0.11
Logistic: C	[1,100]	11.29
Logistic: penalty	L1,L2	L1
Logistic <sub>balanced</sub> : C	[1,100]	14.91
Logistic <sub>balanced</sub> : penalty	L1,L2	L1
Adaboost: No. of estimators	[1,100]	68
Random Forest: depth	[1,5]	3
Random Forest: depth	[1,5]	3

Table 5.4 lists  $F1_{sparse}$  and  $F1_{admit}$  for both the schemes i.e. Simple features as well as Tree features. It should be noted that Logistic regression as used by Water et al [Waters and Miikkulainen, 2013] does

provide the best results if features extracted from the application are fed into the classifier directly, but these decisions can be improved significantly by other forms of modeling such as decision stumps, and even further by using latent variable based generative modeling as detailed in Section 5.4.4.

Figure 5.1 plots  $F1_{sparse}$  as a function of the presence of sparser label. As expected, performance of the model gets better as the sparse label ratio gets closer to 0.5 (i.e. sparse label percentage gets closer to 50%). It is also evident from the figure that all the individual  $F1^{(i)}$ s for decision tree featured model are performing better than simple feature model. The curves, fitted to the individual points by using moving average over a window of 5, show similar trend. Similar is the case for  $F1_{admit}$  when it is plotted against *Acceptance Ratio*, in Figure 5.2.

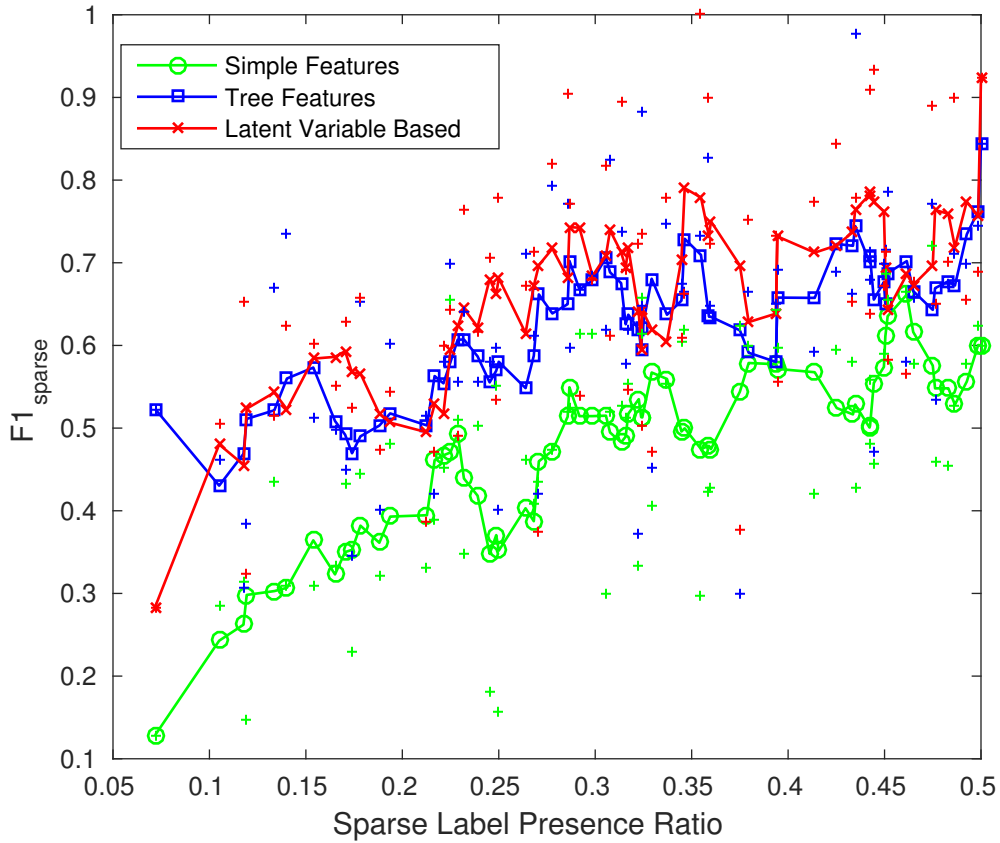


Figure 5.1:  $F1_{sparse}$  as a function of presence of sparse label. The curve is fitted with moving average function of window size 5.

Table 5.4 shows that  $F1_{admit} > F1_{sparse}$  for all the schemes. This is because *Acceptance Ratio* for some universities is greater than 0.5. Predicting *Admit* in this case results in higher  $F1$ . Similar to the axis of Figure 5.1 and Figure 5.2 can be other sorting criteria which can reveal interesting patterns. One such criteria is

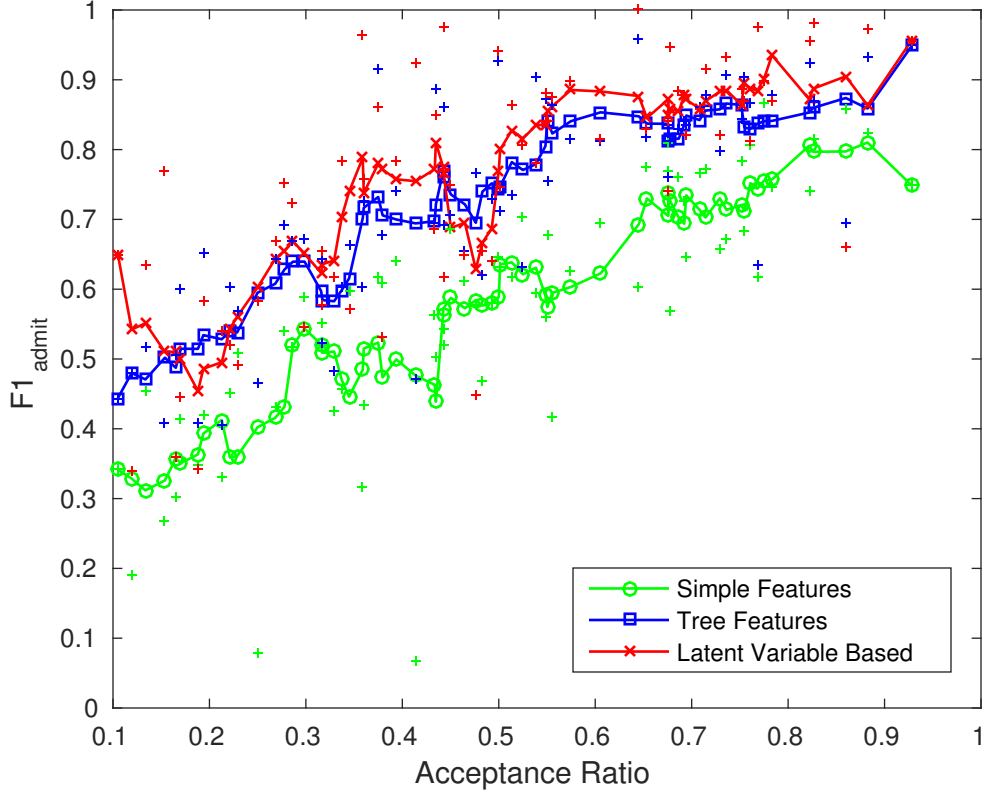


Figure 5.2:  $F1_{admit}$  as a function of Acceptance Ratio. The curve is fitted with moving average function of window size 5.

a proxy of the *reputation* of the university i.e. University ranking. We chose this to be the US News Graduate School Ranking [US News, 2015], primarily, because of its popularity among applicants. Whenever a rank was not available in this resource, a similar resource was consulted [QS Quacquarelli Symonds Limited, 2015, Shanghai, 2015, Webometrics, 2015]. In US News ranks, sometimes multiple adjacently ranked universities were stacked up on a single rank, and the resulting emptied out slots were left vacant. We flattened each stack to its nearest available slots. A high rank means high numerical value of the university rank, which means lower reputation for the university, and vice versa. Thus, a university which is regarded as the best will have the lowest possible rank under this scheme. Interestingly, this ranking scheme has a high positive value of Pearson’s correlation coefficient with *Acceptance Ratio*.

$$\rho(\text{Acceptance Ratio}, \text{University Rank}) = 0.65$$

Such a high correlation suggests a similar trend for  $F1_{admit}$  for University ranking, which is confirmed in Figure 5.3. Figure 5.3 also suggests that it gets easier to predict *Admit* as we go down the rankings.

Table 5.4: F1 over different classifiers and schemes

Classifier / Scheme	Simple Features		Tree Features		Latent Variable Based	
	$F1_{sparse}$	$F1_{admit}$	$F1_{sparse}$	$F1_{admit}$	$F1_{sparse}$	$F1_{admit}$
SVM (Linear Kernel)	0.34	0.52	0.60	<b>0.73</b>	0.66	0.75
SVM <sub>balanced</sub> (Linear Kernel)	0.46	<b>0.58</b>	0.61	0.70	0.66	0.75
SVM (RBF Kernel)	0.18	0.46	0.57	0.70	0.60	0.76
SVM <sub>balanced</sub> (RBF Kernel)	0.09	0.41	0.35	0.57	0.37	0.62
Logistic	0.36	0.56	0.57	0.70	0.60	0.72
Logistic <sub>balanced</sub>	<b>0.48</b>	0.57	0.57	0.68	0.62	0.72
AdaBoost	0.37	0.54	0.55	0.70	0.59	0.71
Random Forest	0.24	0.49	0.59	0.71	0.62	0.74
Random Forest <sub>balanced</sub>	0.42	0.56	<b>0.62</b>	0.72	<b>0.65</b>	<b>0.77</b>

Table 5.5: Discriminative Power of each feature

Index	Feature Name	$F1_{admit}$	Individual F1
1	GPA	0.65	0.65
2	+GRE Quant	0.68	0.53
3	+GRE Verbal	0.71	0.58
4	+GRE AWA	0.73	0.45
5	+TOEFL	0.75	0.58
6	+Program	0.76	0.31
7	+Term	0.76	0.35
8	+Previous Department	0.76	0.42

Despite disagreeing on absolute values, both  $F1_{sparse}$  and  $F1_{admit}$  show similar trends, and given the student’s perspective as explained earlier, we’ll be using  $F1_{admit}$  for the rest of our experiments as well as recommendation system.

### 5.4.3 Feature Selection

These experiments were aimed at understanding the value in each feature. We trained multiple classifiers using single features and evaluated their performances. Then we iteratively added more features to each of the classifiers and evaluated gain in performance. Each feature, when considered individually, is the only classifying parameter. We call its corresponding  $F1_{admit}$  result the Discriminative Power of feature. Although the addition of features one-by-one can theoretically lead to combinatorial explosion, the limited number of original features available in our case prevents this from happening. It was observed that undergraduate GPA has the highest discriminative power and has an average  $F1_{admit}$  over all universities close to 0.65. Figure 5.4 shows how overall  $F1_{admit}$  increases if we sort the features into descending order of discriminative power and keep on adding them to the classifiers.

We observed from experiments that GPA results in a higher  $F1_{admit}$  score compared to any other single

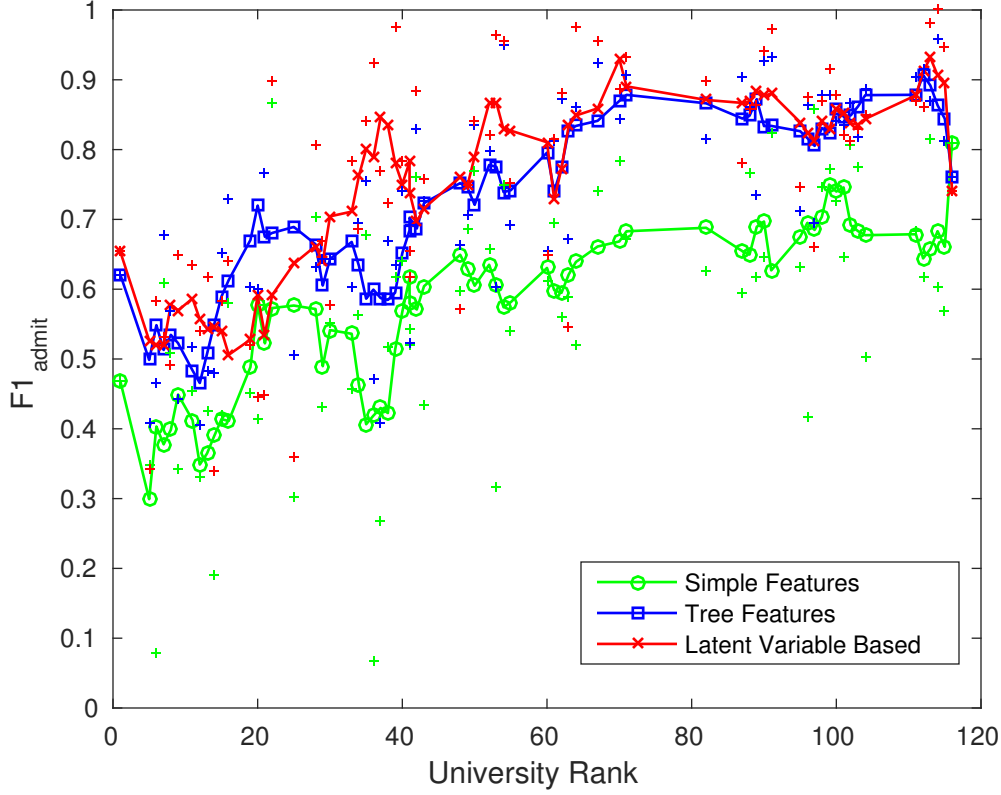


Figure 5.3:  $F1_{admit}$  as a function of Graduate University US News Rank. The curve is fitted with moving average function of window size 5.

feature for any university. Also, if we calculate  $F1_{admit}$  by excluding individual features, exclusion of GPA causes maximum loss in  $F1_{admit}$ . Both of these observations lead to the conclusion that GPA has the highest discriminative power among all available features. The result is intuitive as it validates the expectation that, broadly speaking, GPA is the prime factor in the admission process. It can also be seen that as the number of features that are considered is increased, performance goes up and is the highest when all the features are considered. Table 5.5 lists individual discriminative powers of each feature, as well as cumulative power for  $i$  features i.e. when features  $1, \dots, i$  are used for classification. Table 5.6 lists loss in  $F1_{admit}$  due to exclusion of each feature during classification. Fig 5.4 plots the cumulative discriminative power when we keep on adding features.

#### 5.4.4 Latent Variable Based Approach

The EM model formulated in Section 5.3 was used to cluster students into different groups, representing potential programs. Subsequently, individual classifiers were learned for each of the clusters. Results for



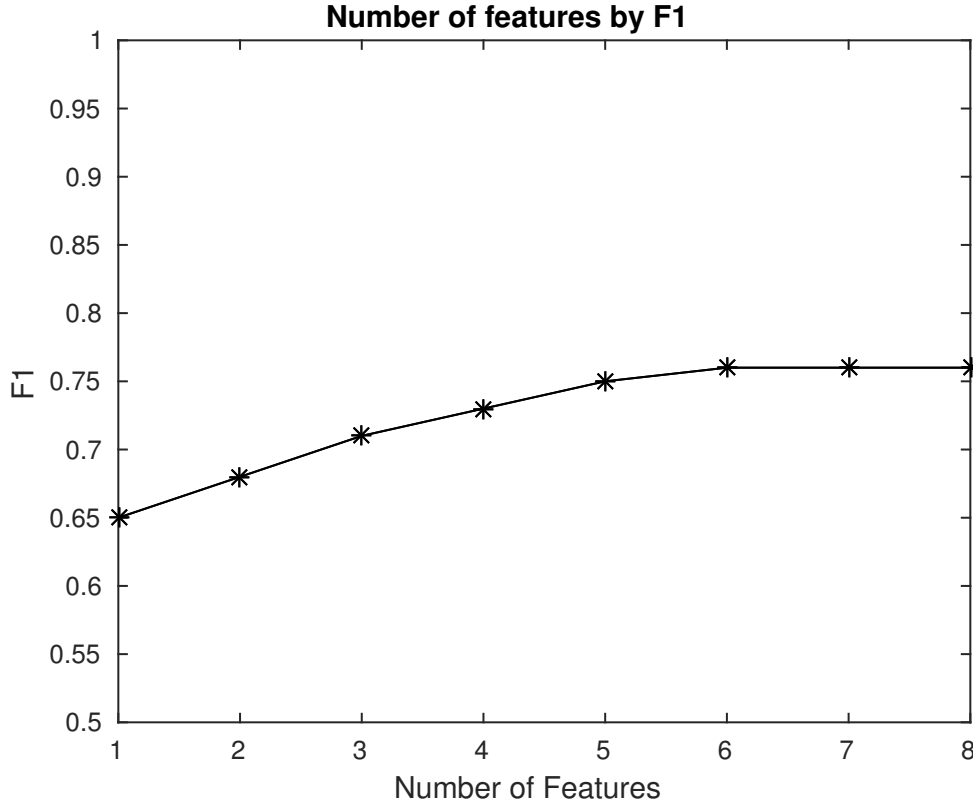


Figure 5.4:  $F1_{admit}$  as we keep on increasing features on top of GPA. Refer to Table 5.5 for feature corresponding to index number on X-axis.

few of the universities for which significant growth was observed are listed in Table 5.7. As per the model assumption, EM bifurcates the data into two clusters ( $z=2$ ), each of which can be separated in a better way than the earlier cumulative cluster thereby increasing the performance of the models significantly. Improvement in  $F1_{admit}$  due to EM clustering is reported in Figure 5.1 and Figure 5.2.

Our model before the use of EM relied on the fact that data for each university is coming from a single source. The improvements in  $F1_{admit}$  as result of splitting the data according to EM formulation indicates that our model is able to capture underlying different distributions of source data. Also, since we know that UIUC offers different degree programs (Professional, Thesis), and CMU offers different specifications (Machine Learning, HCI etc) for the data reported as CS, it is probable that several other universities have more than one underlying distribution because of other factors. Fig 5.3 shows overall increase in  $F1_{admit}$  over all of the universities using EM splitting in two clusters ( $z=2$ ). Table 5.7 reports some of the universities where EM modeling caused roughly 10% or more relative improvement. Although the value of  $z$  is also a tunable parameter of the model, we experimented only with  $z = 2$ . As  $z$ , and correspondingly the number of clusters, increases we expect sparser clusters and hence prone to overfitting.

Table 5.6: F1 without each feature. Less F1 due to missing feature indicates more discriminative power of that feature.

Ignored Feature Name	F1
GPA	0.7234
GRE Quant	0.7415
GRE Verbal	0.7422
GRE AWA	0.7491
TOEFL	0.7455
Program	0.7654
Term	0.7647
Previous Department	0.7518

Table 5.7: Gain in F1 score due to EM clustering

University	Tree	EM + Tree	EM Gain
UCSC	0.58	0.84	0.26
SJSU	0.62	0.86	0.24
UCLA	0.45	0.68	0.23
UMD	0.43	0.66	0.22
SUNY Bingham	0.76	0.94	0.17
UT Austin	0.57	0.74	0.17
UC Boulder	0.80	0.94	0.14
TAMU	0.75	0.86	0.11
UIUC	0.57	0.65	0.08
CMU	0.66	0.71	0.05

### 5.4.5 Understanding Institution Rankings

There is evidence in literature that if an applicant belongs to a *reputed* institution it is regarded as grounds for, at least, some liberality in the way the admission decision would swing [Raghunathan, 2010]. We aimed to qualify this notion by formally validating the assumption - *Undergrad university ranking plays an important role in admissions*. Proving or disproving such an assumption required two-fold experiments:

1. Does university ranking play any role in admission?
2. If yes, what is this rank list?

First, we investigated if the notion of an undergraduate institution’s rank or category even exists. If it does, providing this extra knowledge should help improve the classifier’s performance. Our dataset has applicant records from thousands of undergraduate institutions across many countries. The largest resource for university rankings we were able to find was *Ranking Web of Universities (webometrics)* which ranks 11,701 universities across more than 150 countries and territories [Webometrics, 2015]. Even such a large resource was not exhaustive enough and did not have any ranking information for undergraduate universities of more than 45% applicants. Hence, apart from feeding *webometrics* ranking to

the classifier as feature, we set up several instances of this experiment by using other ranking proxy signals. These proxies included rankings provided by US News [US News, 2015], QS (Top Universities) rankings [QS Quacquarelli Symonds Limited, 2015], Shanghai rankings [Shanghai, 2015] and other lists provided by various private or government agencies such as 'List of Institute of National Importance in India' [MHRD India, 2015]. Since extending the large ranking list of webometrics at such its original finegrain scale was impractical, we divided universities into four categories, as follows:

- *Rank-A*: Institutions ranked as top tier and widely recognized.
- *Rank-B*: Institutions ranked in the middle tier or recognized regionally.
- *Rank-C*: Institutions ranked in the low tier.
- *Rank-D*: Institutions that are neither recognized nor ranked.

Our hypothesis was that reducing the finegrain rankings to broad categories could reduce variance, but such a scheme could easily assign a category value to each university, thus removing the missing value problem for ranking. We expected that if a largely agreed upon rank-list existed, and if our proxies were representative of such a list, then this rank-list should be able to provide gain to the classifier. At the same time, any other list which deviates drastically from such a list should not provide comparable gain during classification.

This category distribution was referred to as '**Original Rank List**' (**ORL**). ORL had following category distribution: Rank-A=47, Rank-B=217, Rank-C=363, Rank-D=2354. Next, we **consciously shuffled** this list using following rules:

1. A university can have either the same category as it originally had, or it can move to its closest category, e.g. Rank-B can move to either Rank-A or Rank-C. The probability of an institution moving to neighbor category is linearly proportional to the target size.
2. Each category still has the same number of institutions as it originally had.

Since, we shuffled the institution categories based on precise rules, we called the result as '**Consciously Shuffled Rank List**' (**CSRL**). We were taking into account that ranking of a university varies with ranking agencies or regions. In addition, we maintained the original category distribution (size of category) of the institutions. '**Randomly Shuffled Same Distribution Rank List**' (**RSSDRL**) was created by assigning a randomly chosen category to each institution but by maintaining original category distribution. Finally, we created a '**Randomly Shuffled Uniform Distribution Rank List**' (**RSUDRL**) by assigning a random

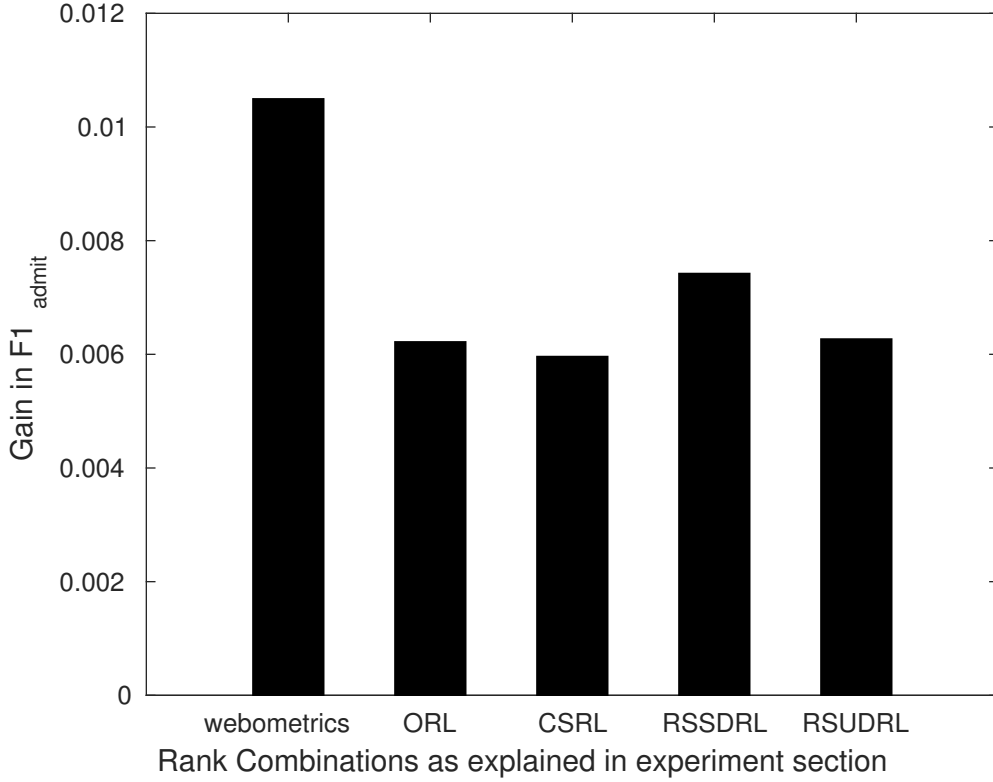


Figure 5.5: Gain in F1 due to various rank-lists

category to each institution, without the constraint of maintaining original distribution. In RSUDRL, each category has uniform probability of occurring within the rank list. These rank-lists will also be released as part of the dataset.

Figure 5.5 indicates that the addition of rank of the undergraduate institution feature led to the gain in performance. We evaluated the gain in terms of statistical significance over 100 iterations and it was significant with  $p - value < 0.0001$ . This leads us to the conclusion that institution rank does play a role in the admission decision. Yet another interesting observation is that there is a comparable gain in webometrics rankings and all of the ways that we created and shuffled the rank-lists, consciously as well as randomly.

Table 5.8 provides gain in  $F1_{admit}$  corresponding to each rank-list for some of the universities. We see that the universities show three types of behavior:

1. There is loss in  $F1_{admit}$  for each of the rank-lists e.g. ASU or CMU
2. Some rank-lists provide gain while others cause loss in  $F1_{admit}$  e.g. Purdue or Brown universities.
3. There is gain in  $F1_{admit}$  for each of the rank-lists e.g. UT Austin or Virginia Tech

Table 5.8: Gain in  $F1_{admit}$  due to various rank-lists (Average over 100 iterations) (In the order of magnitude of  $10^{-3}$ )

University	ORL	CSRL	RSSDRL	RSUDRL
ASU	-0.0047	-0.1093	-0.1307	-0.2356
CMU	-0.0207	-0.7284	-0.0679	-0.2669
Brown University	0.3856	-0.8298	4.0007	-0.0333
Purdue	-2.4818	-1.4320	-1.2061	0.5790
UT Austin	5.7343	4.1302	3.5169	5.2184
Virginia Tech	4.0004	5.4014	5.7771	4.5617

A probable reason for this behavior (1) is that either the universities such as ASU or CMU don't use ranking system at all, and hence providing rank-lists causes the classifier to learn on irrelevant features causing a net loss, or none of the rank-lists is even close to the rank-lists used by these universities. For behavior (2), some of the rank-lists are close to the ones used by these universities while other rank-lists deviate drastically. For observation (3), all of the rank-lists are partially matching to the rank-lists used by these universities. Also, from 3<sup>rd</sup> and 4<sup>th</sup> rows in Table 5.8, it can be seen that one rank-lists provides gain to a specific university, while the other does the same for some other university. Hence, it can be claimed that although the universities use some form of rank-list, there is no consensus over what this actual list is.

#### 5.4.6 Impact of Change in Application Year

In this experiment, we asked the question - Do universities change taste of students over time? Hence, we explored the change in decision to an application in a different application year. Some assume that since there is an increment in the number of applications every year, admissions become more competitive over time. We performed a carefully controlled experiment to test the validity of this hypothesis. In this setting, for every university:

1. Choose a training set (80%) and test set (20%), by random selection.
2. For each record in test set, record the application year, admission decision and prediction of classifier.
3. For each record in test set, change the application year (choose randomly between 2001 and 2015), and record the new prediction, using the classifier used in the previous step.
4. Perform this experiment for  $n(=100)$  iterations.

Our hypothesis was that if yearly factors do not have an effect then, changing application year should not change the decision.

The experiment is aimed at testing the hypothesis that a change in the application year leads to a change in the decision. If an application was correctly classified initially then if there isn't any change in the competition of the application pool then it should still be correctly classified. Over 100 iterations, approximately 60K records were tested, out of which 55K were classified correctly. It was observed that out of these, a vast majority ( $> 98\%$ ) retained the same label even after the change of application year.

We define competition per application to increase if increase in application year changed the decision from *Admit* to *Reject*, and vice versa. Out of all of the decision changes, a record was assigned '+1' if it showed that the competition increased, and '-1' if the competition decreased. Overall sum of these scores for most of the universities, individually, was very close to 0. Whereas, for all of the universities put together, the net sum was -56, which means  $\frac{56}{55,000}$ , approximately 0.1% decrease in competition. Hence, it can be said that decision for an application depends solely on the university and the application, and not the application year.

#### 5.4.7 Which Universities Go Together

One of the unique features of our dataset construction is that applicant records capture various university combinations that the users apply to along with their results. This allowed us to find patterns, and formulate similarities among universities. Apriori algorithm [Agrawal et al., 1994] produces interesting results that are reported in Table 5.10. But Apriori favors heavily populated universities over the less frequent ones. Hence, we expanded our experiments to include null-invariant measures. We computed similarity of two universities based on candidate acceptance using several null-invariant measures such as: AllConf, Jaccard, Cosine Similarity, Kulczynski coefficient, MaxConf defined in [Han et al., 2012]. Results of the experiments are reported in Table 5.9 and 5.10.

As discussed in Chapter 1, applicants have the tendency to apply to universities in buckets (*Dream*, *Reach*, *Safety*). This leads to the hypothesis that since applicants apply in buckets, it should be apparent through the similarity scores of the universities. As the results show this is infact the case. Kulczynski coefficient represents the average of conditional probability conditioned on each of the variables. Table 5.9 lists a few interesting associations based on the kulczynski coefficient and Table 5.10 lists such results for Apriori algorithm.

Table 5.9: Interesting similar universities based on Kulczynski score

University 1	University 2	Kulc
UChicago	CSU	0.286
UNCC	UN Vegas	0.303
CalTech	UCR	0.521
URI	UWisc	0.508

Table 5.10: Universities that go together based on Apriori algorithm

University 1	University 2	Support
UM Twin	SUNY Stony	218
UIC	IU Bloomington	112
Cornell University	SUNY Stony	97
SUNY Buffalo	GMU	75

## 5.5 Recommendations

Since, now we have a system that can predict application decisions for a university, we can utilize it to aid students in making informed choices. In Section 5.4.7, we provide evidence that students apply to universities based on their notion of ‘Dream’, ‘Reachable’ and ‘Safe’ buckets. We include this notion into our algorithm to generate recommendations. Since the classifier system we have is not 100% accurate, it can generate erroneous recommendations if we simply classify for each university and return the results. Fig 5.3 shows that although the trend in university ranks is not strictly monotonous, it becomes very smooth if we cluster neighboring universities and then plot it. Hence, we cluster universities and employ multi-level classification to produce robust results while generating recommendations.

The first level of decision is coarse and the next level result is fine-grained. While classifying coarsely over a range of universities, we mix the records of all universities inside a cluster and train a single classifier on all of them. If the universities in the cluster are similar, the classifier learns the common patterns of admission versus rejection and provides a more general decision than any of the component universities. While in Fine-grained classification an individual classifier is trained for each university.

This algorithm consists of 5 steps:

### 1. University clustering

- Cluster similar universities together based on US News rankings e.g. Universities in rank [1,10] fall into cluster 1, universities from [11,20] fall into cluster 2 and so on.

### 2. Coarse classification

- Using coarse decision, we find the cluster that offers *Admit* and is closest to the top-tier universities. We call this cluster as ‘Reachable’ because it is the best ranked university cluster that can offer Admission.

### 3. Reachable Universities

- Perform fine-grained classification on each of the universities in ‘Reachable’, and return those which produce an *Admit* with highest probability.

### 4. Safe Universities

- We call the cluster next to ‘Reachable’ as ‘Safe’ because it also offers admission and does so with higher probability. Then fine-grained classification is applied on ‘Safe’ to report Safe universities.

### 5. Dream Universities

- For ‘Dream’ universities, we find those universities which are similar to the ones produced by ‘Reachable’ and ‘Safe’ but are towards the top tier universities and hence do not offer admission. These similar universities are based on the higher similarity score based on common admissions.

In step 1, the benefit of using US News rankings is that as the universities get closer to top rank, probability of admission of any candidate decreases. We also verified the same observation from data. As a future work, there many clustering schemes can be employed here, including university similarity scores reported in section 5.4.7, as long as proximity of various clusters is known.



## Chapter 6

# University Perspective

This work was done when the author of this work was an employee of the Department of Computer Science at University of Illinois at Urbana-Champaign, and will not be part of this thesis.

## Chapter 7

# Analysis and Discussion

We started this research with the aim of solving key problems for each of the involved entities in the admission process. For applicants, the problem is not only to get an *Admit*, but also to get an *Admit* from the best possible university. And for universities, the problem is to offer *Admits* to the best applicants while keeping in mind their constraints.

While trying to find answers to these questions, we stumbled upon several interesting findings. We confirmed some of the well-established knowledge while realized that some other assumptions about the process do not stand the empirical tests. The way students currently tackle the problem of uncertainty is through the category assignment of *Dream*, *Reach* and *Safety*. While the assignment has no statistical backup, it is also based on the assumption that the admission behaviour strictly follows the university rankings published by various agencies such as [US News, 2015], [QS Quacquarelli Symonds Limited, 2015] or [Shanghai, 2015]. We found through the experiments in Chapter 5 that the notion of university ranking, although existing, varies from university to university. This revelation contradicts the linear nature of university ranks published by any ranking system or organization. Also, as is suggested by Posselt in [Posselt, 2016], a university might not offer *Admit* to a student if it feels that the student is overqualified and has little chance of actually joining the university. This indicates that what applicants consider *Safety* might not really be *Safety*. All of these revelations tangle the university selection further for the applicant. But we think that we’ve solved this question, at least partially if not fully, for the student. Learning a model for each university with the relevant data provides a good estimate of the student’s chances of admission.

But when looking at the problem from other side of the table, there are different questions. Posselt tells an incident in her book where a student Denpa was rejected at first. But then Denpa’s letter writer called a professor Herald at the university. Herald, in turn, called Vivek, admissions chair, for a meeting. After further discussion and pressure from other professors, eventually Denpa was offered an *Admit*. Posselt also mentions about various biases and reservations of the admission committee and the faculty members serving on the committee. Now, can we hope to perfect this modeling using any of the techniques is a question which begs for a no. During experiments in Chapter 5, we proved that the application year doesn’t affect

the chances of admission for an applicant which means that the admission committee is doing a good job of selection. The implicating claim, though having factual backup, doesn't quantify the term *good* here. To really make such a claim, we need to define a metric for assessing the effectiveness of admission committee. Defining such a metric could also be used for back-propagation of feedback into the admission process itself to overhaul it.

## Chapter 8

# Conclusion and Future Work

This research studies the graduate admission process in American universities using a machine learning approach. Our goal is to build a decision support model that allows candidates to make informed decisions on which schools to apply to, what are their chances of admission, and a slew of other decision-related issues. At the same time, we want to assist the admission committee in universities to make faster and better decisions. Looking at it from student's perspective, we modeled the decision process as a classification problem and presented a system that can achieve high accuracy and can be generalized across multiple universities. By employing many approaches towards solving this problem, such as supervised learning and latent variable based generative modeling, we prove that a mixture of approaches can provide better results than any of the individual approaches. We also provide a dataset that provides avenues for further research. From university's perspective, we show that decision-making can be optimized by using learning techniques to model the task as a ranking problem. We show that significant human effort can be saved by doing so, even without access to all of the application data.

This work can be extended in multiple ways such as towards improving accuracy or validating common notions. Some of the additions of this work may include expanding the EM formulation by modeling further variables such as undergraduate institution ranking mechanism. We proved that every university has a custom ranking mechanism. This mechanism can be modeled as a distribution which can, then, be assigned to one of the hidden variables in the EM model. Theoretically, such a model has more expressive power and can, thus, learn better regarding application decisions. We believe this is but a brisk start to the research that can be performed on the topic. Many more enhancements are possible by expanding the dataset and extracting richer data features from Letters of Recommendation or Statement of Purpose. By asking these questions and providing this dataset, we hope to initiate a discussion that can lead to better understanding of how academia accepts its new members.

# References

- [Agrawal et al., 1994] Agrawal, R., Srikant, R., et al. (1994). Fast algorithms for mining association rules. *Proc. 20th int. conf. very large data bases, VLDB*, 1215:487–499.
- [Breiman, 2001] Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- [Breiman et al., 1984] Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). Classification and regression trees. wadsworth. *Belmont, CA*.
- [Bruggink and Gambhir, 1996] Bruggink, T. H. and Gambhir, V. (1996). Statistical models for college admission and enrollment: A case study for a selective liberal arts college. *Research in Higher Education*, 37(2):221–240.
- [Chang and Lin, 2011] Chang, C.-C. and Lin, C.-J. (2011). Libsvm: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):27.
- [Dempster et al., 1977] Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38.
- [Efron and Tibshirani, 1986] Efron, B. and Tibshirani, R. (1986). Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statist. Sci.*, 1(1):54–75.
- [Freund and Schapire, 1995] Freund, Y. and Schapire, R. E. (1995). A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational learning theory*, pages 23–37. Springer.
- [Grove and Roth, 2001] Grove, A. and Roth, D. (2001). Linear concepts and hidden variables. *Machine Learning*, 42(1/2):123–141.
- [Han et al., 2012] Han, J., Kamber, M., and Pei, J. (2012). 6 - mining frequent patterns, associations, and correlations: Basic concepts and methods. In Kamber, J. H. and Pei, J., editors, *Data Mining (Third Edition)*, The Morgan Kaufmann Series in Data Management Systems, pages 243 – 278. Morgan Kaufmann, Boston, third edition edition.
- [Kassegne, 2016] Kassegne, S. (2016). Edulix - premier site for scholars - 'education crowdsourced'. *Web*.
- [Krishnamoorthy, 2013] Krishnamoorthy, S. (2013). Acceptance rates, apparently, are poor predictors of getting in. *The New York Times*.
- [MHRD India, 2015] MHRD India (2015). List of Institutes of National Importance in India. *Web*.
- [Moore, 1998] Moore, J. S. (1998). An expert system approach to graduate school admission decisions and academic performance prediction. *Omega*, 26(5):659 – 670.
- [Murthy and Salzberg, 1995] Murthy, K. V. S. and Salzberg, S. L. (1995). *On growing better decision trees from data*. PhD thesis, Citeseer.
- [National Science Foundation, 2014] National Science Foundation (2014). Survey of Graduate Students and Postdoctorates in Science and Engineering. *Web*.

- [Pedregosa et al., 2011] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- [Posselt, 2016] Posselt, J. (2016). *Inside graduate admissions: Merit, diversity, and faculty gatekeeping*. Harvard University Press.
- [QS Quacquarelli Symonds Limited, 2015] QS Quacquarelli Symonds Limited (2015). Top Universities — Worldwide university rankings, guides and events. *Web*.
- [Quinlan, 1986] Quinlan, J. R. (1986). Induction of decision trees. *Machine learning*, 1(1):81–106.
- [Raghunathan, 2010] Raghunathan, K. (2010). Demystifying the American Graduate Admissions Process. *Web*.
- [Shanghai, 2015] Shanghai (2015). Academic Ranking of World Universities. *Web*.
- [Smola and Schölkopf, 2004] Smola, A. J. and Schölkopf, B. (2004). A tutorial on support vector regression.
- [The Grad Cafe, 2015] The Grad Cafe (2015). The Grad Cafe. *Web*.
- [US News, 2015] US News (2015). Education Rankings and Advice. *Web*.
- [Waters and Miikkulainen, 2013] Waters, A. and Miikkulainen, R. (2013). Grade: Machine learning support for graduate admissions. In *Proceedings of the 25th Conference on Innovative Applications of Artificial Intelligence*.
- [Webometrics, 2015] Webometrics (2015). Ranking Web of Universities. *Web*.